

## A APPENDIX

### A.1 ADDITIONAL FIGURES

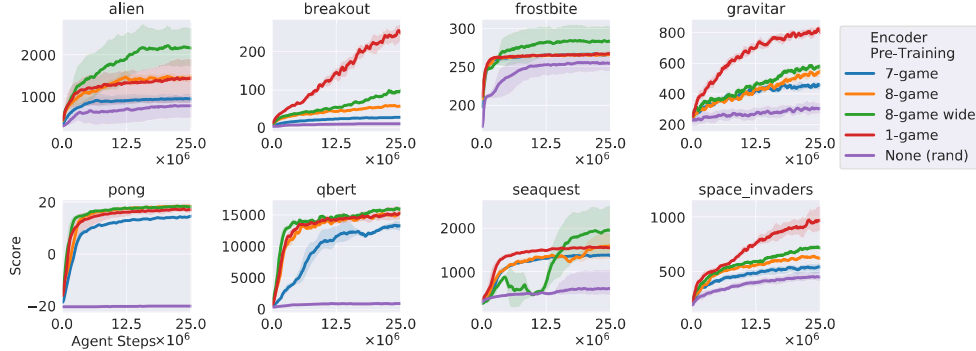


Figure 11: RL using multi-task encoders (all with weights frozen) for eight Atari games gives mixed performance, partially improved by increased network capacity (8-game-wide). Training on 7 games and testing on the held-out one yields diminished but non-zero performance, showing some limited feature transfer between games.

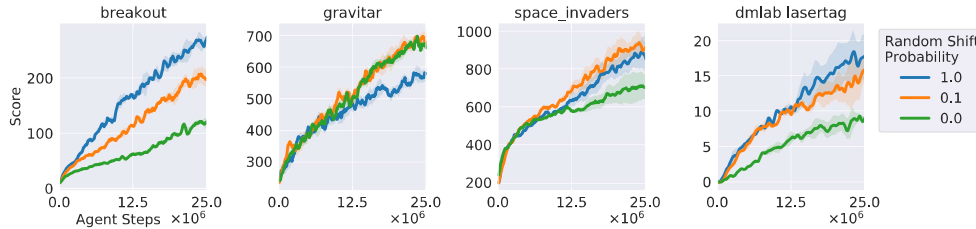


Figure 12: Random shift augmentation helps in some Atari games and hurts in others, but applying with probability 0.1 is a performant middle ground. DMLab benefits from random shift. (Offline pre-training.)

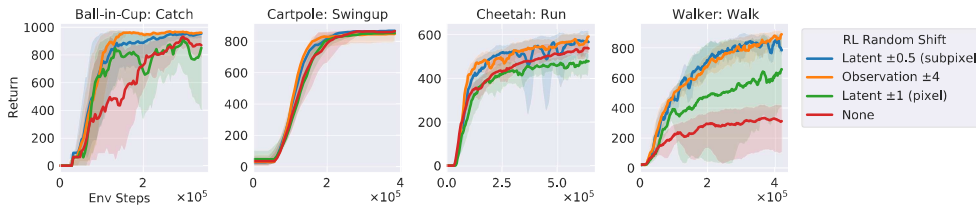


Figure 13: Even after pre-training encoders for DMControl using random shift, RL requires augmentation—our subpixel augmentation acts on the (compressed) latent image, permitting its use in the replay buffer.

In subpixel random shift, new pixels are a linearly weighted average of the four nearest pixels to a randomly chosen coordinate location. We used uniformly random horizontal and vertical shifts, and tested maximum displacements in  $(\pm) \{0.1, 0.25, 0.5, 0.75, 1.0\}$  pixels (with “edge” mode padding  $\pm 1$ ). We found 0.5 to work well in all tested domains, restoring the performance of raw image augmentation but eliminating convolutions entirely from the RL training updates.

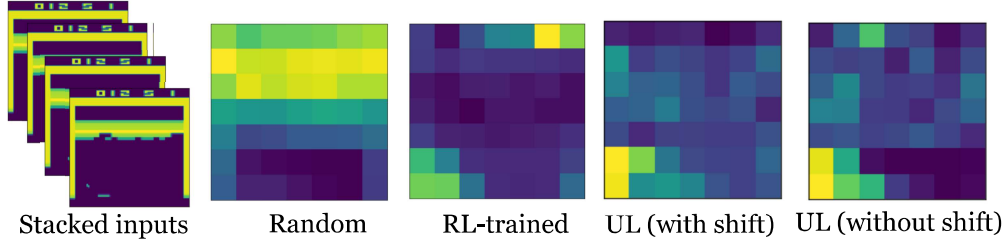


Figure 14: Attention map in BREAKOUT which shows the RL-trained encoder focusing on game score, whereas UL encoder focuses properly on the paddle and ball.

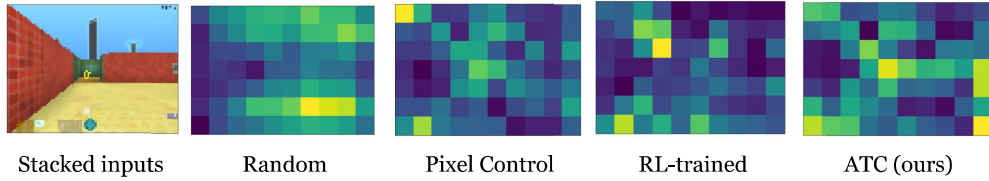


Figure 15: Attention map in LASERTAG. UL encoder with pixel control focuses on the score, while UL encoder with the proposed ATC focuses properly on the coin similar to RL-trained encoder.

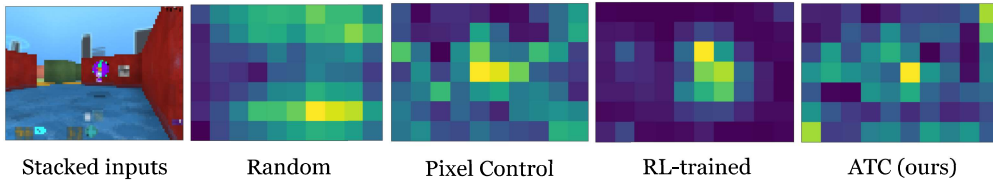


Figure 16: Attention map in the LASERTAG which shows that UL encoders focus properly on the enemy similar to RL-trained encoder.

## A.2 RL SETTINGS

Table 1: DMControl, RAD-SAC Hyperparameters.

HYPERPARAMETER	VALUE
OBSERVATION RENDERING	(84, 84), RGB
RANDOM SHIFT PAD	$\pm 4$
REPLAY BUFFER SIZE	1e5
INITIAL STEPS	1e4
STACKED FRAMES	3
ACTION REPEAT	2 (FINGER, WALKER) 8 (CARTPOLE) 4 (REST)
OPTIMIZER	ADAM
$(\beta_1, \beta_2) \rightarrow (f_\theta, \pi_\psi, Q_\phi)$	(.9, .999)
$(\beta_1, \beta_2) \rightarrow (\alpha)$	(.5, .999)
LEARNING RATE $(f_\theta, \pi_\psi, Q_\phi)$	2e-4 (CHEETAH) 1e-3 (REST)
LEARNING RATE $(\alpha)$	1e-4
BATCH SIZE	512 (CHEETAH, PENDULUM) 256 (REST)
$Q$ FUNCTION EMA $\tau$	0.01
CRITIC TARGET UPDATE FREQ	2
CONVOLUTION FILTERS	[32, 32, 32, 32]
CONVOLUTION STRIDES	[2, 2, 2, 1]
CONVOLUTION FILTER SIZE	3
ENCODER EMA $\tau$	0.05
LATENT DIMENSION	50
HIDDEN UNITS (MLP)	[1024, 1024]
DISCOUNT $\gamma$	.99
INITIAL TEMPERATURE	0.1

Table 2: Atari, PPO Hyperparameters.

HYPERPARAMETER	VALUE
OBSERVATION RENDERING	(84, 84), GREY
STACKED FRAMES	4
ACTION REPEAT	4
OPTIMIZER	ADAM
LEARNING RATE	2.5e-4
PARALLEL ENVIRONMENTS	16
SAMPLING INTERVAL	128
LIKELIHOOD RATIO CLIP, $\epsilon$	0.1
PPO EPOCHS	4
PPO MINIBATCHES	4
CONVOLUTION FILTERS	[32, 64, 64]
CONVOLUTION FILTER SIZES	[8, 4, 3]
CONVOLUTION STRIDES	[4, 2, 1]
HIDDEN UNITS (MLP)	[512]
DISCOUNT $\gamma$	.99
GENERALIZED ADVANTAGE ESTIMATION $\lambda$	0.95
LEARNING RATE ANNEALING	LINEAR
ENTROPY BONUS COEFFICIENT	0.01
EPISODIC LIVES	FALSE
REPEAT ACTION PROBABILITY	0.25
REWARD CLIPPING	$\pm 1$
VALUE LOSS COEFFICIENT	1.0

Table 3: DMLab, PPO Hyperparameters.

HYPERPARAMETER	VALUE
OBSERVATION RENDERING	(72, 96), RGB
STACKED FRAMES	1
ACTION REPEAT	4
OPTIMIZER	ADAM
LEARNING RATE	$2.5e-4$
PARALLEL ENVIRONMENTS	16
SAMPLING INTERVAL	128
LIKELIHOOD RATIO CLIP, $\epsilon$	0.1
PPO EPOCHS	1
PPO MINIBATCHES	2
CONVOLUTION FILTERS	[32, 64, 64, 64]
CONVOLUTION FILTER SIZES	[8, 4, 3, 3]
CONVOLUTION STRIDES	[4, 2, 1, 1]
HIDDEN UNITS (LSTM)	[256]
SKIP CONNECTIONS	CONV 3, 4; LSTM
DISCOUNT $\gamma$	.99
GENERALIZED ADVANTAGE ESTIMATION $\lambda$	0.97
LEARNING RATE ANNEALING	NONE
ENTROPY BONUS COEFFICIENT	0.01 (EXPLORE)
	0.0003 (LASERTAG)
VALUE LOSS COEFFICIENT	0.5

## A.3 ONLINE UL SETTINGS

Table 4: Common ATC Hyperparameters.

HYPERPARAMETER	VALUE
RANDOM SHIFT PAD	$\pm 4$
LEARNING RATE	$1e-3$
LEARNING RATE ANNEALING	COSINE
TARGET UPDATE INTERVAL	1
TARGET UPDATE $\tau$	0.01
PREDICTOR HIDDEN SIZES, $h_\psi$	[512]

Table 5: DMControl ATC Hyperparameters.

HYPERPARAMETER	VALUE
RANDOM SHIFT PROBABILITY	1
BATCH SIZE	AS RL (INDIVIDUAL OBSERVATIONS)
TEMPORAL SHIFT, $k$	1
MIN AGENT STEPS TO UL	$1e4$
MIN AGENT STEPS TO RL	$1e4$
UL UPDATE SCHEDULE	AS RL (2X CHEETAH)
LATENT SIZE	128

Table 6: Atari ATC Hyperparameters.

HYPERPARAMETER	VALUE
RANDOM SHIFT PROBABILITY	0.1
BATCH SIZE	512 (32 TRAJECTORIES OF 16 TIME STEPS)
TEMPORAL SHIFT, $k$	3
MIN AGENT STEPS TO UL	$5e4$
MIN AGENT STEPS TO RL	$1e5$
UL UPDATE SCHEDULE	ANNEALED QUADRATICALLY FROM 6 PER SAMPLER ITERATION ( $1e4$ ONCE AT $1e5$ STEPS FOR WEIGHT INITIALIZATION)
LATENT SIZE	256

Table 7: DMLab ATC Hyperparameters.

HYPERPARAMETER	VALUE
RANDOM SHIFT PROBABILITY	1
BATCH SIZE	512 (INDIVIDUAL OBSERVATIONS)
TEMPORAL SHIFT, $k$	3
MIN AGENT STEPS TO UL	$5e4$
MIN AGENT STEPS TO RL	$1e5$
UL UPDATE SCHEDULE	2 PER SAMPLER ITERATION
LATENT SIZE	256

#### A.4 OFFLINE PRE-TRAINING DETAILS

We conducted coarse hyperparameter sweeps to tune each competing UL algorithm. In all cases, the best setting is the one shown in our comparisons.

When our VAEs include a time difference between input and reconstruction observations, we include one hidden layer with action additionally input between the encoder and decoder. We tried both 1.0 and 0.1 KL-divergence weight in the VAE loss, and found 0.1 to perform better in both DMControl and Atari.

**DMControl** For the VAE, we experimented with 0 and 1 time step difference between input and reconstruction target observations and training for either  $1e4$  or  $5e4$  updates. The best settings were 1-step temporal, and  $5e4$  updates, with batch size 128. ATC used 1-step temporal,  $5e4$  updates (although this can be significantly decreased), and batch size 256 (including CHEETAH). The pre-training data set consisted of the first  $5e4$  transitions from a RAD-SAC agent learning each task, including  $5e3$  random actions. Within this span, CARPOLE and BALL-IN-CUP learned completely, but WALKER and CHEETAH reached average returns of 514 and 630, respectively (collected without the compressive convolution).

**DMLab** For Pixel Control, we used the settings from Hessel et al. (2019) (see the appendix therein), except we used only empirical returns, computed offline (without bootstrapping). For CPC, we tried training batch shapes,  $batch \times time$  in  $(64, 8)$ ,  $(32, 16)$ ,  $(16, 32)$ , and found the setting with rollouts of length 16 to be best. We contrasted all elements of the batch against each other, rather than only forward contrasts. In all cases we also used 16 steps to warmup the LSTM. For all algorithms we tried learning rates  $3e-4$  and  $1e-3$  and both  $5e4$  and  $1.5e5$  updates. For ATC and CPC, the lower learning rate and higher number of updates helped in LASERTAG especially. The pretraining data was  $125e3$  samples from partially trained RL agents receiving average returns of 127 and 6 in EXPLORE\_GOAL\_LOCATIONS\_SMALL and LASERTAG\_THREE\_OPPONENTS\_SMALL, respectively.

**Atari** For the VAE, we experimented with 0, 1, and 3 time step difference between input and reconstruction target, and found 3 to work best. For ST-DIM we experimented with 1, 3, and 4 time steps differences, and batch sizes from 64 to 256, learning rates  $1e-3$  and  $5e-4$ . Likewise, 3-step delay worked best. For the inverse model, we tried 1- and 3-step predictions, with 1-step working better overall, and found random shift augmentation to help. For pixel control, we used the settings in Jaderberg et al. (2017), again with full empirical returns. We ran each algorithm for up to  $1e5$  updates, although final ATC results used  $5e4$  updates. We ran each RL agent with and without observation normalization on the latent image and observed no difference in performance. Pretraining data was  $125e3$  samples sourced from the replay buffer of DQN agents trained for  $15e6$  steps with epsilon-greedy  $\epsilon = 0.1$ . Evaluation scores were:

Table 8: Atari Pre-Training Data Source Agents.

GAME	EVALUATION SCORE
ALIEN	1,800
BREAKOUT	279
FROSTBITE	1,400
GRAVITAR	390
PONG	18
QBERT	8,800
SEAQUEST	11,000
SPACE INVADERS	1,200